

WHAT IS CLAIMED IS:

1 1. A method for eliminating indistinguishable differentials
2 from a direct comparison of a pair of data matrices comprising a plurality
3 of data points, wherein each data point in the first matrix of the pair has a
4 corresponding data point in the second matrix of the pair, the method
5 comprising:

6 (a) ranking the data points of each matrix from highest to
7 lowest according to intensity such that a plot of rank versus intensity of
8 the data points for each matrix provides an experimental curve for each
9 matrix;

10 (b) fitting a smooth curve to each experimental curve to
11 provide a model curve for each matrix, each model curve comprising a
12 first section separated from a second section by an inflection point;

13 (c) eliminating any pair of corresponding data points for
14 which each data point of the pair is below the inflection point of its model
15 curve; and

16 (d) eliminating any pair of corresponding data points for
17 which the rank of each data point of the pair is below a selected cutoff
18 rank between r_{\min} and r_{\max} , where r_{\min} is the rank of the data point at the
19 minimum of the derivative of one of the model curves and r_{\max} is the rank
20 of the highest ranking data point on the model curve.

1 2. The method of claim 1, wherein each model curve has
2 an analytical rank determined by the equation:

3
$$f'(\text{analytical rank}) = C ((f(r_{\max})/r_{\max})),$$

4 where $f'(\text{analytical rank})$ is the value of the first derivative of
5 the model curve at the analytical rank, r_{\max} is the rank of a data point
6 having a rank higher than CR on the model curve, and $f(r_{\max})$ is the value
7 of the model curve at r_{\max} ;

8 and further wherein the cutoff rank is the larger analytical
9 rank of the two model curves.

1 3. The method of claim 1, wherein r_{\max} is the rank of a
2 data point having an intensity in the upper 30% of the model curve.

1 4. The method of claim 1, wherein r_{\max} is the rank of the
2 highest ranking data point on the model curve.

1 5. The method of claim 1, further comprising normalizing
2 the model curves by transforming the intensity values for the data points
3 in one of the model curves to become equal to the intensity values of the
4 data points of the same rank in the other model curve.

1 6. The method of claim 5, wherein the intensity values
2 for the model curve corresponding to the less noisy data matrix are
3 transformed to become equal to the intensity values of model curve
4 corresponding to the noisier data matrix.

1 7. The method of claim 1, wherein the model curves are
2 generated by fitting the experimental curves with the equation:

$$3 \quad f(x) = \left(\frac{a_1}{x + a_2} + \frac{x}{r_{\max} - x + a_3} - a_4 \right) * a_5 * \left(\frac{1}{1 + \left(\frac{a_7}{x} \right)^{a_6}} + \frac{a_8}{1 + \left(\frac{a_{10}}{x} \right)^{a_9}} - \frac{a_{11}}{x + a_{12}} \right) * \left(1 + \frac{a_{13}}{1 + \left| 1 - \frac{a_{15}}{x} \right|^{a_{14}}} \right) + \left(\frac{1}{1 + \left(\frac{a_{17}}{x} \right)^{a_{16}}} - a_{18} \right) * a_{19}$$

4 where x and $f(x)$ refer to the rank and logarithmic value of the intensity
5 at rank x , respectively, (a_1, \dots, a_{19}) are variables and r_{\max} is the rank of the
6 highest ranking data point in the data matrix.

1 8. The method of claim 2, wherein C has a value from
2 about 0.3 to about 0.4.

1 9. The method of claim 2, wherein C has a value from
2 about 0.35 to about 0.37.

1 10. The method of claim 1, further comprising measuring a
2 background level for each data points in each matrix and eliminating any
3 pair of corresponding data point for which at least one data point of the
4 pair has a background level that lies outside two or more standard
5 deviations from the mean value of the background measurements for its
6 matrix.

1 11. The method of claim 1, wherein each data matrix is
2 generated from an individual sample.

1 12. The method of claim 1, wherein each data matrix is
2 generated from a composite sample comprising no more than about 100
3 individual samples.

1 13. The method of claim 1, wherein the data matrices are
2 generated from biological samples and the data points represent the
3 expression levels of biomolecules produced by the biological samples.

1 14. The method of claim 1, wherein the biological samples
2 are selected from the group consisting of cell samples, tissue samples,
3 biological fluid samples, DNA samples and RNA samples.

1 15. The method of claim 1, wherein the data matrices are
2 generated from a gene expression profiling experiment and the data points
3 represent gene expression levels.

1 16. The method of claim 1, wherein the data matrices are
2 generated from a protein expression profiling experiment and the data
3 points represent protein expression levels.

1 17. The method of claim 1, wherein the data matrices are
2 generated from a nuclei acid sequence expression profiling system and
3 the data points represent oligonucleotide expression levels.

1 18. A method for eliminating false differentials from a
2 direct comparison of three or more replicate data matrix pairs, each data
3 matrix comprising a plurality of data points, wherein each data point in
4 the first matrix of the pair has a corresponding data point in the second
5 matrix of the pair and each pair of corresponding data points in each
6 matrix pair has a corresponding pair of corresponding data points in every
7 other matrix pair, the method comprising:

8 (a) eliminating indistinguishable differentials from the data
9 matrices according to the method of claim 1;

10 (b) determining an intensity ratio for each remaining pair
11 of corresponding data points in each data matrix pair, wherein each ratio
12 is categorized as less than one, greater than one, or equal to one; and

13 (c) eliminating any corresponding pairs of corresponding
14 data points for which the intensity ratios for the corresponding pairs fall
15 into the same category for less than one half of the corresponding pairs.

1 19. The method of claim 18, comprising eliminating
2 corresponding pairs of corresponding data points if the intensity ratios for
3 the corresponding pairs fall into the same category for less than 75
4 percent of the corresponding pairs.

1 20. The method of claim 18, comprising eliminating
2 corresponding pairs of corresponding data points if the intensity ratios for
3 the corresponding pairs fall into the same category for less than 100
4 percent of the corresponding pairs.

1 21. The method of claim 18, further comprising eliminating
2 any data points whose $\log_2(\text{intensity ratio})$ values lie outside the largest
3 standard deviation of all the remaining $\log_2(\text{intensity ratio})$ values
4 multiplied by a constant.

1 22. The method of claim 21, wherein the constant is at
2 least about 2.

1 23. The method of claim 18, wherein the intensity ratios
2 for the remaining data points comprise no more than 1 % false
3 differentials.

1 24. The method of claim 18, wherein the intensity ratios
2 for the remaining data points comprise no more than 0.1 % false
3 differentials.

1 25. The method of claim 18, wherein the intensity ratios
2 for the remaining data points comprise no more than 0.01 % false
3 differentials.

1 26. The method of claim 18, wherein the data matrices
2 are generated from a gene expression profiling experiment and the
3 intensity ratios represent gene expression ratios.

1 27. The method of claim 26, wherein the gene expression
2 ratios remaining after the method is applied provide information about
3 gene function behind a biological phenotype.

1 28. The method of claim 26, wherein the gene expression
2 ratios remaining after the method is applied provide information about the
3 genetic networks behind a biological phenotype.

1 29. A method for measuring the quality of a direct
2 comparison of a pair of data matrices comprising a plurality of data
3 points, wherein each data point in the first matrix of the pair has a
4 corresponding data point in the second matrix of the pair, the method
5 comprising:

6 (a) eliminating indistinguishable differentials from the data
7 matrices according to the method of claim 1; and

8 (b) calculating a noise factor (NF) for the pair of matrices
9 according to the equation:

$$10 \quad NF = \sqrt{\frac{\sum_{i=1}^n (r_{1i} - r_{2i})^2}{n}} * \frac{K}{(r_{\max} - CR)}$$

11 wherein CR is the cutoff rank, n is the total number of
12 remaining data points whose ranks are larger than CR in both arrays, r_{\max}
13 is the rank of a data point having a rank higher than CR, r_{1i} is the rank of
14 data point i in the first data matrix of the pair, r_{2i} is the rank of the
15 corresponding data point in the second data matrix of the pair, and K is a
16 constant.

1 30. The method of claim 29, wherein r_{\max} is the rank of a
2 data point having an intensity in the upper 30% of the model curve.

1 31. The method of claim 29, wherein r_{\max} is the rank of the
2 highest ranking data point on the model curve.

1 32. The method of claim 29, wherein the data matrices
2 are generated from biological samples and the data points represent the
3 expression levels of biomolecules produced by the biological samples.

1 33. The method of claim 29, wherein the data matrices
2 are generated from a gene expression profiling experiment and the
3 intensity ratios represent gene expression ratios.

1 34. The method of claim 33, wherein the gene expression
2 ratios remaining after the method is applied provide information about
3 gene function behind a biological phenotype.

1 35. The method of claim 33, wherein the gene expression
2 ratios remaining after the method is applied provide information about the
3 genetic networks behind a biological phenotype.